

Lecture 24

(1)

(b) Non equal covariance metrics $\Sigma_1 \neq \Sigma_2$

(1) Quadratic rule based on MECM

- Population case

Allocate X_0 to

$$\begin{aligned} \pi_1 & \text{ if } -\frac{1}{2} X_0' (\Sigma_1^{-1} - \Sigma_2^{-1}) X_0 + (\mu_1' \Sigma_1^{-1} - \mu_2' \Sigma_2^{-1}) X_0 - K \geq \ln Y \\ \pi_2 & \text{ if } \dots \dots \dots < \ln Y \end{aligned}$$

$$\text{Where } K = \frac{1}{2} \ln \frac{|\Sigma_1|}{|\Sigma_2|} + \frac{1}{2} (\mu_1' \Sigma_1^{-1} \mu_1 - \mu_2' \Sigma_2^{-1} \mu_2)$$

- Sample case

1°) $\hat{\mu}_i = \bar{X}_i, \hat{\Sigma}_i = S_i$

2°) Allocate X_0 to

$$\begin{aligned} \pi_1 & \text{ if } -\frac{1}{2} X_0' (S_1^{-1} - S_2^{-1}) X_0 + (\bar{X}_1' S_1^{-1} - \bar{X}_2' S_2^{-1}) X_0 - \hat{K} \geq \ln Y \\ & \text{ if } \dots \dots \dots < \ln Y \end{aligned}$$

$$\text{Where } K = \frac{1}{2} \ln \frac{|S_1|}{|S_2|} + \frac{1}{2} (\bar{X}_1' S_1^{-1} \bar{X}_1 - \bar{X}_2' S_2^{-1} \bar{X}_2)$$

(2) Practical issues

- check normality first
- If necessary, employ transformation method to get more normal data
- test equality of $\Sigma_1 = \Sigma_2$ (Chapter 6.6) Likelihood ratio test (chi-square test)
- Try different classification procedures and compare their performance and select the best.

11.4 Evaluating classification procedures

(a) Optimal error rate (misclassification probability)

(1) Total probability of misclassification (TPM)

$$\begin{aligned} \text{TPM} &= P_1 P(2|1) + P_2 P(1|2) \\ &= P_1 \int_{R_2} f_1(x) dx + \int_{R_1} P_2 f_2(x) dx \\ &= P_1 + \int_{R_1} P_2 f_2(x) - P_1 f_1(x) dx \end{aligned}$$

(2) Optimal error rate (OER)

• To minimize TPM we should select

$$R_1 : \{x : P_2 f_2(x) - P_1 f_1(x) \leq 0\}$$

$$\Leftrightarrow \{x : \frac{f_1(x)}{f_2(x)} \geq \frac{P_2}{P_1}\}$$

$$\begin{aligned} \bullet \text{ OER} &= P_1 + \int_{\{x : \frac{f_1(x)}{f_2(x)} \geq \frac{P_2}{P_1}\}} (P_2 f_2(x) - P_1 f_1(x)) dx \\ &\quad \uparrow \\ &\quad \text{minimized TPM} \end{aligned}$$

(3) Actual error rate (AER)

• In practice, after we have \hat{R}_1 and \hat{R}_2

We obtain an estimate of OER, called actual error rate (AER)

$$\text{AER} = P_1 \int_{\hat{R}_2} f_1(x) dx + P_2 \int_{\hat{R}_1} f_2(x) dx$$

(3)

- AER still depends on $f_i(x)$, p_i , $i=1, 2$

(b) Apparent error rate (APER)

(1) APER

- Applicable for any classification rules
- It is the fraction of observations in the "training" sample that are misclassified
- Confusion matrix \swarrow predicted membership

		π_1	π_2	
Actual membership \rightarrow	π_1	n_{1c}	n_{1m}	n_1
	π_2	n_{2m}	n_{2c}	n_2

$$\bullet \text{ APER} = \frac{n_{1m} + n_{2m}}{n_1 + n_2}$$

(2) Advantages and disadvantages

- Easy to calculate (popular)
- But usually can underestimate the AER (unless n_1 and n_2 are very large)
- Reuse of data (building and judging)

(c) Some other procedure and expected actual error rate

• Splitting sample

1°) Split sample into

$\left\{ \begin{array}{l} \text{training sample} \rightarrow \text{to construct classification rule} \\ \text{validation sample} \rightarrow \text{to check performance} \\ \text{(test)} \end{array} \right.$

2°) Disadvantages

requires large sample

may lose valuable information

• Cross-validation to estimate expected AER

(1) cross-validation (jack knife, leave-one-out, hold out)

1°) start with π_1 , hold out one point from π_1

develop a classification rule based on $n_1 - 1, n_2$

2°) classify the "hold out" point using the rule in 1°)

3°) repeat 1° and 2°) until all points in π_1 are classified

Let $n_{1m}^{(H)}$ be the # of misclassified points in π_1

4°) repeat 1° - 3°) for π_2 , obtain $n_{2m}^{(H)}$

(2) Estimate of expected AER

$$\hat{E}(\text{AER}) = \frac{n_{1m}^{(H)} + n_{2m}^{(H)}}{n_1 + n_2}$$

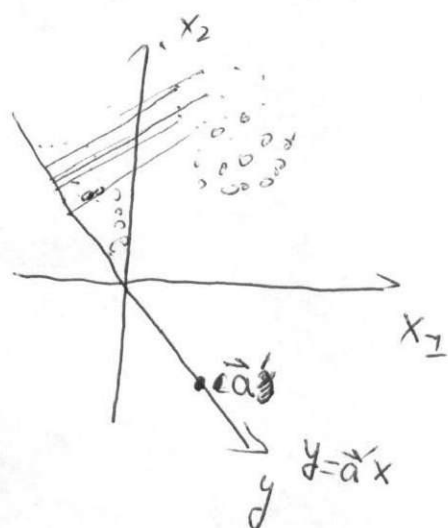
(5)

11.5 Fisher's Discriminant function - Separation of two populations ($\Sigma_1 = \Sigma_2$) (No normal assumption)

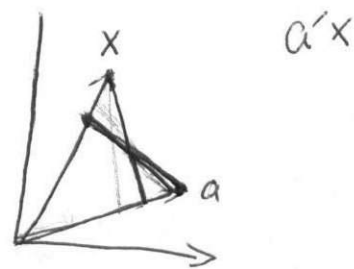
(a) basic idea and classification rule

(1) Multivariate to univariate via transformation

- Transform the multivariate observation x from π_1 and π_2 to univariate y such that the y 's derived from π_1 and π_2 are separated as far as possible



$a'x$



$$\vec{a} = (a_1, a_2)'$$

$$x = (x_1, x_2)'$$

1°) Let x_{11}, \dots, x_{1n_1} and x_{21}, \dots, x_{2n_2} are samples from π_1 and π_2

Let a be a vector, project x_{ij} to a , $i=1,2$, $j=1, \dots, n$

Let $y = a'x$ then

$$y_{11} = a'x_{11} \quad y_{21} = a'x_{21}$$

$$y_{1n_1} = a'x_{1n_1} \quad y_{2n_2} = a'x_{2n_2}$$

π_1

π_2

$$\bar{y}_1 = a'\bar{x}_1$$

$$\bar{y}_2 = a'\bar{x}_2$$

(6)

$$2^o) \text{ Let } S_{2y}^2 = \frac{1}{n_2-1} \sum_{j=1}^{n_2} (y_{2j} - \bar{y}_2)^2 = a' S_{2x} a$$

$$S_y^2 = \frac{(n_1-1) S_{1y}^2 + (n_2-1) S_{2y}^2}{n_1+n_2-2} = a' S_{\text{pooled}} a$$

$$S_{\text{pooled}} = \frac{(n_1-1) S_{1x}^2 + (n_2-1) S_{2x}^2}{n_1+n_2-2}$$

Define Separation = $\frac{\bar{y}_1 - \bar{y}_2}{s_y}$

3^o) Find \hat{a} that maximizes

$$\text{Separation}^2 = \frac{(\bar{y}_1 - \bar{y}_2)^2}{s_y^2} = \frac{(a'(\bar{x}_1 - \bar{x}_2))^2}{a' S_{\text{pooled}} a}$$

4^o Since $\frac{(a'(\bar{x}_1 - \bar{x}_2))^2}{a' S_{\text{pooled}} a}$

$$= \frac{((S_{\text{pooled}}^{-\frac{1}{2}} a)' S_{\text{pooled}}^{-\frac{1}{2}} (\bar{x}_1 - \bar{x}_2))^2}{(S_{\text{pooled}}^{-\frac{1}{2}} a)' (S_{\text{pooled}}^{-\frac{1}{2}} a)}$$

$$= \frac{(u'v)^2}{\|u\|^2} \leq \frac{(\|u\| \|v\|)^2}{\|u\|^2} = \|v\|^2$$

$$= (\bar{x}_1 - \bar{x}_2)' S_{\text{pooled}}^{-\frac{1}{2}} S_{\text{pooled}}^{-\frac{1}{2}} (\bar{x}_1 - \bar{x}_2)$$

The equality is satisfied only if

$$u = cv \quad \text{Let } c=1$$

$$\Rightarrow S_{\text{pooled}}^{-\frac{1}{2}} a = S_{\text{pooled}}^{-\frac{1}{2}} (\bar{x}_1 - \bar{x}_2)$$

$$\Rightarrow \hat{a} = S_{\text{pooled}}^{-1} (\bar{x}_1 - \bar{x}_2)$$

(7)

(2) Fisher's discrimination function (classification) rule

- $\hat{a}'X$ is called Fisher's discriminant function

Let $\hat{y} = \hat{a}'x$, then allocate x_0 to

$$\begin{cases} \pi_1 : \text{if } \hat{y}_0 = \hat{a}'x_0 \geq \frac{1}{2}(\bar{y}_1 + \bar{y}_2) = \hat{a}'(\frac{\bar{x}_1 + \bar{x}_2}{2}) \\ \pi_2 : \text{if } \hat{y}_0 < \frac{1}{2}(\bar{y}_1 + \bar{y}_2) \end{cases}$$

(equivalent to
MECM rule for
 $r=1$, under normal
populations with
 $\Sigma_1 = \Sigma_2 = \Sigma$)

(b) Test for $\mu_1 = \mu_2$ (under normal assumption)

(1) Maximum separation and the test (See chapter 6)

- Consider $H_0: \mu_1 = \mu_2$

If $\pi_i: N_p(\mu_i, \Sigma)$, then

$$\left(\frac{1}{n_1} + \frac{1}{n_2}\right)^{-1} D^2 \sim \frac{(n_1 + n_2 - 2)P}{n_1 + n_2 - P - 1} F_{P, n_1 + n_2 - P - 1}$$

$$D^2 = (\bar{x}_1 - \bar{x}_2)' S_{\text{pooled}}^{-1} (\bar{x}_1 - \bar{x}_2)$$

- reject H_0 if D^2 is large (not empty good classification)

- If separation is not significant - the search for a useful classification rule will probably prove fruitless.

11.6 classification with several population (g populations, $g \geq 2$)

(a) MECM method and rule

<1> the method

- $P_i (i=1, \dots, g)$ prior probability

$$P(K|i) = P(X \text{ allocated to } \pi_K | x \in \pi_i)$$

$$= \int_{R_K} f_i(x) dx \quad i, K=1, \dots, g$$

R_K : the set of x 's classified as π_K

$C(K|i)$ = the cost of allocating x to π_K while it belongs to π_i

- $ECM(1) = \sum_{K=1}^g C(K|1) P(K|1)$

$$= \sum_{K=1}^g C(K|1) P(K|1)$$

$$ECM(i) = \sum_{K \neq i} C(K|i) P(K|i)$$

$$ECM = \sum_{i=1}^g P_i ECM(i)$$

$$= \sum_{i=1}^g P_i \left(\sum_{K \neq i} C(K|i) P(K|i) \right)$$

(2) The rule

• Allocate x to π_k if

$$\frac{\sum_{\substack{i=1 \\ i \neq k}}^g P_i f_i(x) c(k|i)}{\sum_{i=1}^g P_i f_i(x) c(j|i)} \leq \frac{\sum_{\substack{i=1 \\ i \neq j}}^g P_i f_i(x) c(j|i)}{\sum_{i=1}^g P_i f_i(x) c(j|i)} \quad \text{for all } j=1, 2, \dots, g \text{ and } j \neq k$$

• If all misclassification costs are equal, then allocate x to π_k if

$$P_k f_k(x) \geq P_i f_i(x) \quad \text{for all } i \neq k$$

\Leftrightarrow

$$p(\pi_k | x) = \frac{P_k f_k(x)}{\sum_j P_j f_j(x)} \geq \frac{P_i f_i(x)}{\sum_j P_j f_j(x)} = p(\pi_i | x) \quad \text{for all } i \neq k$$

"posterior probability"

\Leftrightarrow

$$\ln P_k f_k(x) \geq \ln P_i f_i(x) \quad \text{for all } i \neq k$$

(b) classify normal populations

\hookrightarrow Unequal Σ_i

• Quadratic discriminant Score

1°) Population Case.

Consider $\pi_i : N_p(\mu_i, \Sigma_i)$, $i=1, \dots, g$, then

$$\ln P_i f_i(x) = \ln P_i + \ln f_i(x)$$

$$= \ln P_i - \frac{p}{2} \ln 2\pi - \frac{1}{2} \ln |\Sigma_i| - \frac{1}{2} (x - \mu_i)' \Sigma_i^{-1} (x - \mu_i)$$

According to the MFCM rule with $c(k|i)$ all the same,

Allocate x to π_k if

$$\ln P_k f_k(x) = \max_i \ln P_i f_i(x)$$

Define the quadratic discriminant score as

$$d_i^Q(x) = -\frac{1}{2} \ln |\Sigma_i| - \frac{1}{2} (x - \mu_i)' \Sigma_i^{-1} (x - \mu_i) + \ln P_i$$

2°) Sample case

$$\hat{d}_i^Q(x) = -\frac{1}{2} \ln |\hat{\Sigma}_i| - \frac{1}{2} (x - \bar{x}_i)' \hat{\Sigma}_i^{-1} (x - \bar{x}_i) + \ln P_i$$

• MFCM rule (with equal misclassification costs)

2°) Sample case

$$x \rightarrow \pi_k \text{ if } \hat{d}_k(x) = \max_{1 \leq i \leq g} \hat{d}_i^Q(x)$$

1°) population case

$$x \rightarrow \pi_k \text{ if } d_k^Q(x) = \max_{1 \leq i \leq g} d_i^Q(x)$$

(2) Equal Σ_i linear

• linear discriminant ~~case~~ score

1°) population case

$$\begin{aligned} \ln P_i f_i(x) &= \ln P_i - \left(\frac{p}{2}\right) \ln 2\pi - \frac{1}{2} \ln |\Sigma| - \frac{1}{2} (x - \mu_i)' \Sigma^{-1} (x - \mu_i) \\ &= \left[\ln P_i - \left(\frac{p}{2}\right) \ln 2\pi \right] - \frac{1}{2} \ln |\Sigma| - \frac{1}{2} x' \Sigma^{-1} x \\ &\quad + \mu_i' \Sigma^{-1} x - \frac{1}{2} \mu_i' \Sigma^{-1} \mu_i + \ln P_i \end{aligned}$$

(17)

Define the linear discriminant score as

$$d_i(x) = (\Sigma^{-1} \mu_i)' x - \frac{1}{2} \mu_i' \Sigma^{-1} \mu_i + \ln P_i, \quad i=1, \dots, g$$

2°) Sample case

$$\hat{d}_i(x) = (S_{\text{pooled}}^{-1} \bar{x})' x - \frac{1}{2} (\bar{x}_i)' S_{\text{pooled}} \bar{x}_i + \ln P_i$$

$$S_{\text{pooled}} = \frac{(n_1-1)S_1 + \dots + (n_g-1)S_g}{(n_1 + \dots + n_g) - g}$$

• MECM with equal $c(k|i)$'s

1°) population case

$$x \rightarrow \pi_k \quad \text{if } d_k(x) = \max_i d_i(x)$$

2°) sample case

$$x \rightarrow \pi_k \quad \text{if } \hat{d}_k(x) = \max_i \hat{d}_i(x)$$

$$\Leftrightarrow \underbrace{-\frac{1}{2} (x - \bar{x}_k)' S_{\text{pooled}}^{-1} (x - \bar{x}_k)}_{= \max_i (-\frac{1}{2} D_i^2 + \ln P_i)} + \ln P_k \quad P_k^2$$

(c) Examples

11.7 Fisher's discriminant functions for separation of several populations

(A) Idea and classification rule

(1) Idea and Assumption

- Idea: Find a few linear combinations

$$a_1'X, a_2'X, \dots$$

Which could be used to represent populations and separate populations as much as possible

- Assumption $\Sigma_1 = \Sigma_2 = \dots = \Sigma_g = \Sigma$

(2) population case

- Let μ_i be the mean of π_i and $\bar{\mu} = \frac{1}{g} \sum_{i=1}^g \mu_i$

- Let $B_{\mu} = \sum_{i=1}^g (\mu_i - \bar{\mu})(\mu_i - \bar{\mu})'$ measures between group variability
between groups sum of cross products

- Consider $Y = \alpha'X$ then

$$\text{Var}(Y) = \alpha' V(X) \alpha = \alpha' \Sigma \alpha = \sigma_Y^2$$

$$E(Y) = \alpha' E(X) = \alpha' \mu_i \quad \text{if } X \in \pi_i \\ = \mu_{iY}$$

$$\text{Let } \bar{\mu}_Y = \frac{1}{g} \sum_{i=1}^g \mu_{iY} = \frac{1}{g} \sum_{i=1}^g \alpha' \mu_i = \alpha' \bar{\mu}$$

• Define
$$\text{Separation}^2 = \frac{\sum_{i=1}^g (\mu_i - \bar{\mu})^2}{\sigma^2} \leftarrow \text{sum of squared distances from populations to overall mean of } Y$$

$$= \frac{a' B a}{a' \Sigma a} \leftarrow \text{The common variability within groups}$$

• To find a to maximize separation²

(3) Sample case

• Let x_{i1}, \dots, x_{in_i} be a sample from π_i

Let $\bar{x}_i = \frac{1}{n_i} \sum_{j=1}^{n_i} x_{ij} \rightarrow \text{estimate } \mu_i$

$\bar{x} = \frac{\sum_{i=1}^g n_i \bar{x}_i}{\sum_{i=1}^g n_i} \rightarrow \text{estimate } \bar{\mu}$

$$\text{Spooled} = \frac{\sum_{i=1}^g (n_i - 1) S_i}{\sum_{i=1}^g n_i - g} = \frac{W}{\sum_{i=1}^g n_i - g} \rightarrow \text{estimate } \Sigma$$

(*)
$$B = \sum_{i=1}^g (n_i) (\bar{x}_i - \bar{x})(\bar{x}_i - \bar{x})' \rightarrow B_{\mu}$$

• Find \hat{a} to minimize

$$\begin{aligned} \frac{a' B a}{a' \text{Spooled } a} &\Leftrightarrow \frac{a' B a}{a' W a} = \frac{(W^{\frac{1}{2}} a)' W^{-\frac{1}{2}} B W^{-\frac{1}{2}} (W^{\frac{1}{2}} a)}{(W^{\frac{1}{2}} a)' (W^{\frac{1}{2}} a)} \\ &= \frac{u' W^{-\frac{1}{2}} B W^{-\frac{1}{2}} u}{u' u} = \frac{u'}{\|u\|} W^{-\frac{1}{2}} B W^{-\frac{1}{2}} \frac{u}{\|u\|} \leq \lambda_{\max}(W^{-\frac{1}{2}} B W^{-\frac{1}{2}}) \\ &= \lambda_{\max}(W^{-1} B) \end{aligned}$$

$$W^{\frac{1}{2}} \hat{a} = e_1 (W^{-\frac{1}{2}} B W^{\frac{1}{2}})$$

$$W^{-\frac{1}{2}} B W^{\frac{1}{2}} \cdot W^{\frac{1}{2}} \hat{a} = \lambda_{\max} \cdot W^{\frac{1}{2}} \hat{a}$$

$$W^{-1} B \hat{a} = \lambda_{\max} \hat{a}$$

$$\Rightarrow \hat{a} = e_{\max} (W^{-1} B)$$

(4) Fisher's sample linear discriminants

• Let $\hat{\lambda}_1 \geq \hat{\lambda}_2 \geq \dots \geq \hat{\lambda}_s > 0$, $s \leq \min(g-1, p)$

be nonzero eigenvalues of $W^{-1}B$ and $\hat{e}_1, \dots, \hat{e}_s$ be the corresponding eigen vectors

Then $\hat{a}_1 = \hat{e}_1$ maximizes

$$\frac{a' B a}{a' W a} \quad \text{or} \quad \frac{a' B a}{a' S_{\text{pooled}} a}$$

• $\hat{y}_1 = \hat{a}_1' x = \hat{e}_1' x$ is called the sample first discriminant
 $\hat{y}_2 = \hat{a}_2' x = \hat{e}_2' x$ ————— second —————

$\hat{y}_k = \hat{a}_k' x = \hat{e}_k' x$ ————— kth —————

(5) Fisher's classification procedure based on sample discriminant

• Let $\hat{y}_j = \hat{a}_j' x$, $\bar{y}_{kj} = \hat{a}_j' \bar{x}_k$ $j=1, \dots, r$, $r \leq S$

1°) Allocate $x \rightarrow \pi_k$ if

$$\sum_{j=1}^r (\hat{y}_j - \bar{y}_{kj})^2 = \sum_{j=1}^r [\hat{a}_j' (x - \bar{x}_k)]^2$$

$$\leq \sum_{j=1}^r (\hat{y}_j - \bar{y}_{ij})^2 = \sum_{j=1}^r [\hat{a}_j' (x - \bar{x}_i)]^2, \quad \forall i \neq k$$

2° Let $g = \begin{pmatrix} \hat{y}_1 \\ \vdots \\ \hat{y}_r \end{pmatrix}$, $\bar{y}_k = \begin{pmatrix} \bar{y}_{k1} \\ \vdots \\ \bar{y}_{kr} \end{pmatrix} = \begin{pmatrix} \hat{a}_1' \bar{x}_k \\ \vdots \\ \hat{a}_r' \bar{x}_k \end{pmatrix}$

$$= \begin{pmatrix} \hat{a}_1' x \\ \vdots \\ \hat{a}_r' x \end{pmatrix}$$

Then allocate $x \rightarrow \pi_k$ if

$$\| \hat{y} - \bar{y}_k \|^2 \leq \| \hat{y} - \bar{y}_i \|^2, \quad \forall i \neq k$$

code11-2

```
options ls=85 ps=65;
title1 'SAS DISCRIM example 2';
title2 h=1 'Admission data for graduate school of business; JW: T11-6';
data admission;
  infile 'Z:\\My Documents\\Teaching\\Stat524\\Fall 2010\\Data set\\T11-6.dat' firstobs=1;
  input gpa gmat status$;
  gpa1=200*gpa;
run;

proc print data=admission;
run;

proc plot data=admission;
  plot gpa*gmat=status;
run;

title2 "scatter plot of gpa and gmat";
%plotit(data=admission,labelvar=_blank_, symvar=status,
plotvars=gpa1 gmat, color=black, colors=blue);
run;

proc discrim data=admission method=normal pool=test
  wcov pcov manova listerr crosslisterr;
  class status;
  var gpa gmat;
  *priors equal/propotional/'1'=0.3 '2'=0.5 '3'=0.1;
run;

proc discrim data=admission method=normal pool=yes wcov pcov
  listerr crosslisterr;
  class status;
  var gpa gmat;
run;

proc discrim data=admission method=normal pool=no wcov pcov
  listerr crosslisterr;
  class status;
  var gpa gmat;
run;
```


Chapter 12 Cluster Analysis

< Introduction

● clustering and classification

- (1) Classification, Gives the number of groups, assigns new observation to one of the groups.
- (2) clustering: (grouping): No assumption on the number of groups or the structure of groups. searches for some 'natural' groups of items (variables) based on similarities (or dissimilarities)

12.2 Similarity measures

(a) two types of problems in cluster analysis

(1) cluster obs (items, units, cases) based on some sort of "distance" (numerical measurements)

(2) cluster variables based on correlation coefficients

(b) Similarity measurements for pairs of items

(1) Euclidean distance

$$X = \begin{pmatrix} x_1 \\ \vdots \\ x_p \end{pmatrix}, Y = \begin{pmatrix} y_1 \\ \vdots \\ y_p \end{pmatrix}$$

$$d(X, Y) = \sqrt{(x_1 - y_1)^2 + \dots + (x_p - y_p)^2} = \left(\sum (x_i - y_i)^2 \right)^{\frac{1}{2}}$$

(2)

(2) L^m distance

$$d(x, y) = \left(\sum_{i=1}^p |x_i - y_i|^m \right)^{\frac{1}{m}} \quad m \geq 1$$

(3) Statistical distance

$$d(x, y) = \sqrt{(x - y)^T S^{-1} (x - y)}$$

but then what's S ?

(4) other distance

- Sometimes no meaningful p-dim measurement exists for items.

One may introduce binary variables based on the presence and absence of a characteristic

1°) gender handedness

$$\text{item } i \quad M(1) \quad R(0) \rightarrow (1, 0) = x_i$$

$$\text{item } k \quad L(0) \quad R(1) \rightarrow (0, 1) = x_k$$

$$x_{ij} = \begin{cases} 1 & \text{if item } i \text{ has } j\text{th character} \\ 0 & \text{otherwise} \end{cases}$$

\nearrow i th item number \nwarrow j th character

2°) dissimilarity

$$d^2(i, k) = d^2(x_i, x_k) = \sum_{j=1}^p (x_{ij} - x_{kj})^2$$

count of mismatches

- disadvantages:

Weighting 1-1 and 0-0 matches equally

• Similarity coefficients

1°) Contingency table

		item K		
		1	0	
item i	1	a	b	a+b
	0	a c	d	c+d
		a+c	b+d	p = a+b+c+d

2°) Similarity coefficients

$d(c, k) = :$

(1) $\frac{a+d}{p} = \frac{a+d}{a+b+c+d}$

(2) $\frac{2(a+d)}{2(a+d)+c+b}$

(3) $\frac{a+d}{a+d+2(c+b)}$

(4) $\frac{a}{p}$

(5) $\frac{a}{a+b+c}$

(6) . . .

(C) similarity measures for pairs of variables

(4)

 X_i, X_k

① Correlation Coefficient

② similarity coefficients for binary variables

• Contingency table

variable i	variable k			n	obs	$X_i \dots X_k$
	1	0				1 0
1	a	b	a+b	}	}	1 0 1 0
0	c	d	c+d			1 0 1 1
	a+c	b+d	n=a+b+c+d			1 0 1 1
						1 0 1 1

• Similarity coefficients

 $d(i, k) :=$

$$① \quad r_{ij} = \frac{\sum (X_{ij} - \bar{X}_i)(X_{kj} - \bar{X}_k)}{\sqrt{\sum (X_{ij} - \bar{X}_i)^2} \sqrt{\sum (X_{kj} - \bar{X}_k)^2}} = \frac{ad - bc}{\sqrt{(a+b)(c+d)(a+c)(b+d)}}$$

$$② \quad \frac{a+d}{n}$$

$$③ \quad \frac{2(a+d)}{2(a+d)+c+b}$$

(d) Matrix D :

(Symmetric)

item (variable)	item (variable)				
	1	2	3	4	...
1	$d(1,1)$	-		$d(1,4)$	
2	$d(2,1)$	$d(2,2)$		$d(2,4)$	-
...					

- Hierarchical clustering methods

- (a) Agglomerative and divisive hierarchical methods

- ↳ agglomerate

Initially, each object is a cluster, then a series of successive mergers

- (2) Divisive

Initially, only one cluster that contains all objects, then a series of successive divisions.

- (b) Agglomerative hierarchical procedures

- ↳ clustering algorithm for grouping N objects.

- Input: A matrix $D = \{d_{ij}\}_{N \times N}$

(distance, similarity, coordinates)

- Output: A tree diagram (dendrogram)

- Steps

- 1° start with N clusters, and a matrix $D = \{d_{ij}\}$

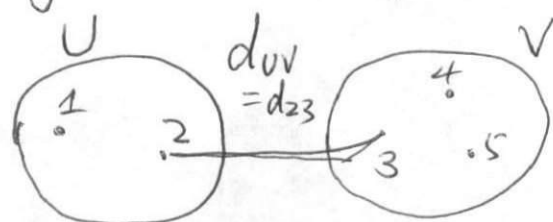
- 2° Search D for the "most similar" (nearest) pair of clusters U and V

- 3° Merge U and V , labeled the new cluster UV , update D

- 4° Repeat 2° and 3° $N-1$ times

- (2) Single linkage (method)

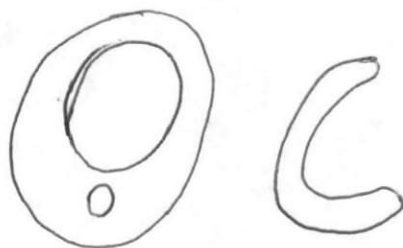
- Distance between two clusters is identified to be the smallest distance
- Merge nearest neighbors (pairs with the smallest distance or largest similarity)



$$d_{UV} = \min_{\substack{u \in U \\ v \in V}} d_{uv}$$

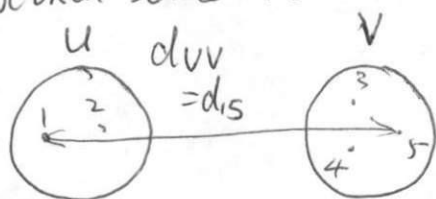
(recognized, distinct, defect)

- Can't discern poorly-separated clusters but good for unelliptical ones



(3) Complete linkage

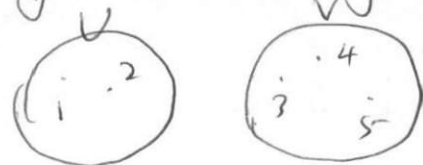
- Distance between two clusters is the maximum distance
- Merge nearest neighbors (so that all items in a cluster are within some maximum distance (or minimum similarity))



$$d_{UV} = \max_{\substack{u \in U \\ v \in V}} d_{uv}$$

(4) Average linkage

- Distance between two clusters is the average distance
- Merge nearest neighbors



$$d_{UV} = \frac{\sum_{i=1}^2 \sum_{j=1}^3 d_{ij}}{6} = \frac{\sum_{u \in U} \sum_{v \in V} d_{uv}}{N_U N_V}$$

(5) Ward's Hierarchical cluster's method

- Matrix $D = \{d_{ij}\}$ is coordinate matrix
variables

	1	2	3	...	p
Items	1 x_{11}	x_{12}	—	—	x_{1p}
2					
:					
:					
N	x_{N1}	—	—	—	x_{Np}

- Error sum of squares (ESS)

$$ESS_K = \sum_{i \in K} (x_i - \bar{x}_K)(x_i - \bar{x}_K) \quad K=1 \sim g$$

of groups

$$ESS = \sum_K ESS_K$$

- Initially, N clusters

$$ESS_K = 0, \quad K=1 \dots N \quad ESS=0$$

- Merge pairs of clusters U and V such that the resulting increasing in ESS is minimum

- Good for observations that are roughly elliptically shaped.

10 Examples

9

<1>

	①	②	③	④	⑤
D =	0	0	0	0	0
	9	0	0	0	0
	3	7	0	0	0
	6	5	9	0	0
	11	10	2	8	0

⑧ Single linkage

• merge ③ and ⑤ → (3,5), ①, ②, ④

• update

	①	②	④
(3,5)	0	0	0

⑩ Single linkage

10) Single linkage

①	3	0	0
②	7	9	0
④	8	6	5

→ (1,3,5), ②, ④

2° complete linkage

	①	②	④
(3,5)	0	0	0
①	11	0	0
②	10	9	0
④	9	6	5

→ (3,5), (2,4), ①

3° Average linkage

	①	②	④
(3,5)	0	0	0
①	7	0	0
②	8.5	9	0
④	8.5	6	5

→ (3,5), (2,4), ①

12.4 Nonhierarchical clustering method

(a) objective and advantage

<1> objective

- cluster items (not variabbs) into K disjoint groups
- K is specified in advance or determined during clustering process

<2> advantage

- Good for large data set

SAS cluster Fastclus

(b) K-means method

<1> algorithm

- Input: A coordinate matrix and a number K
- Output: K (disjoint) clusters
- Step:

1°) partition items into K initial clusters.

calculate the centroid (mean) for each cluster

2°) Reassign an item to the cluster whose centroid (mean) is nearest to the item.

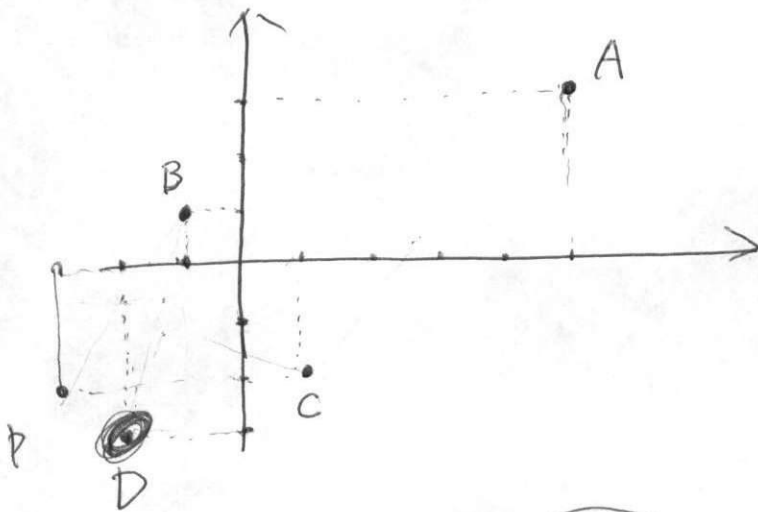
3°) Repeat 2°) untill no more reassignment take place

(2) Example

Variables

Items		x_1	x_2
A		5	3
B		-1	1
C		1	-2
D		-3	-2

- objective: group items into 2 clusters



- start with (AB) and (CD)
centroid A B C D $(\text{distance})^2$

$$(AB) \quad (2, 2) \quad 10 \quad 10 \quad 17 \quad 41$$

$$(CD) \quad (-1, -2) \quad 6^2 + 5^2 = 61 \quad 9 \quad 4 \quad 4$$

- end up ~~with~~ assigning B to $(CD) \Rightarrow (A) \quad (BCD)$

$$\begin{array}{ccccc} \text{centroid} & A & B & C & D \\ A & 5 & 3 & 0 & 6^2 + 2^2 = 40 & 41 & 89 \end{array}$$

$$(BCD) \quad (-1, -1) \quad 6^2 + 4^2 = 52 \quad 4 \quad 5 \quad 5$$

- no more assigning, end up with $(A) \quad (BCD)$

• However, if we start with (AC) and (BD) (11)
end up with (AC) and (BD)

• start with B (ACD) \rightarrow (AC) (BD)

(C) Comments

(1) choice of K (no clear cut)

- Subject-matter knowledge
- data-based appraisals
- specify
- try different K 's

(2) For incommensurable responses
standardize observations

(3) K-mean method is not robust against outliers

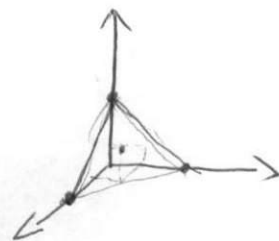
- should try different distance measures

L^m -distance $C \propto m < 1$

- use multivariate median as centroid

$$\arg \min \sum \|x_i - \theta\|$$

- Use multivariate weighted mean



(1,0,0)
(0,1,0) (0,0,1)